Word  Embeddings.

a  glass  of  orange  _____ ?

$e_1$
$e_2$
$e_3$
$e_4$

10000 × 1.
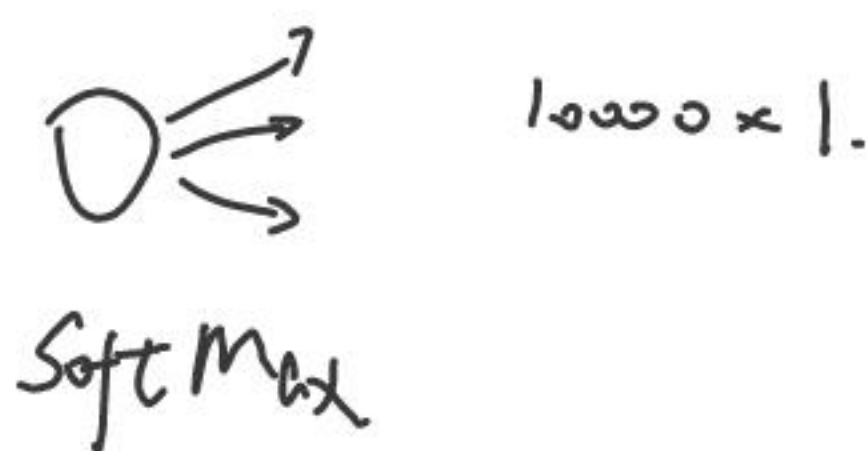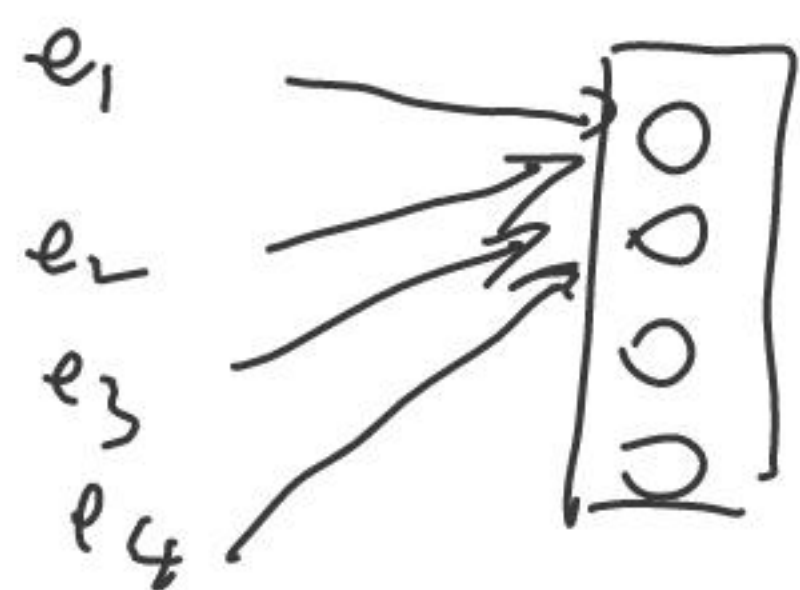
SoftMax

windows  size = 4          vocabulary  size = 10,000.

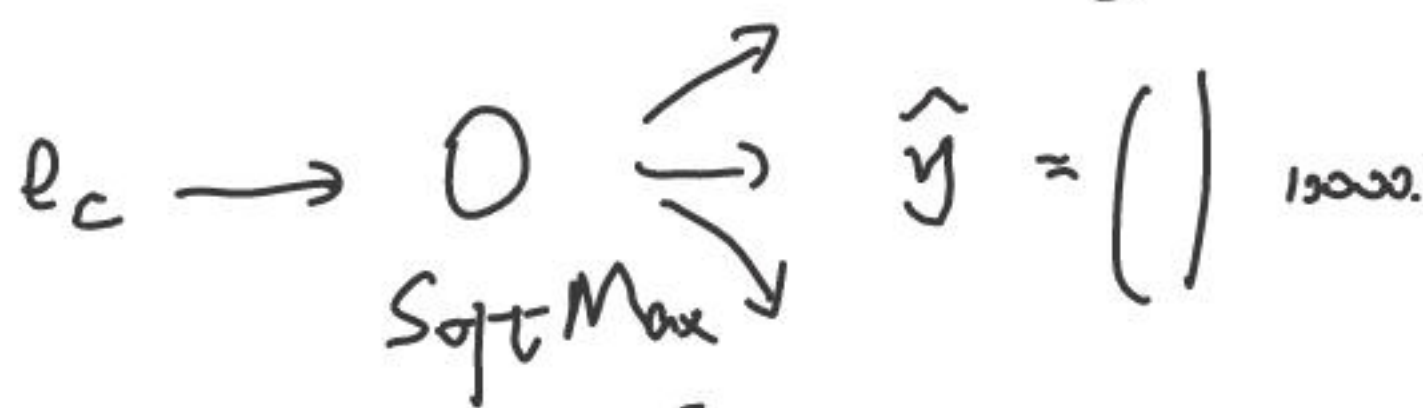训练 E, W, b.

2003,  neural  probabilistic  language  model.

Skip - Gram.

给定中心词 预测上下文 . "context"  "target".

a    glass    of    orange    juice.

$e_c \longrightarrow O \quad \hat{y} = \begin{pmatrix} \\ \end{pmatrix}$ 10000.

SoftMax

$$P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_j e^{\theta_j^T e_c}}.$$

Problems: 计算成本高.  ⇒ Hierachi  SoftMax

对 c 的采样也用启发式,        树型分类器.

避免 " the 、of、a、an".
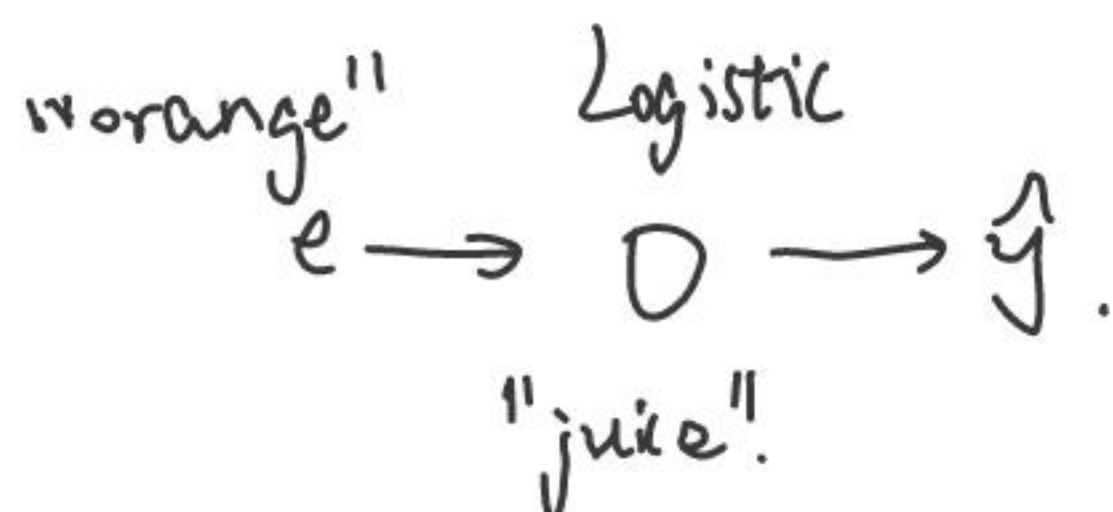
# Negative Sampling.

正样本：负样本 = 1：k, 随机采样.

$$P(y=1 \mid t, c) = \sigma(\theta_t^T \cdot e_c).$$

"orange"     Logistic

$$e \longrightarrow \bigcirc \longrightarrow \hat{y}.$$

"juice"

采样策略：在均匀分布与频率分布之间最"合适".

## Sentiment Classification.

1):

the        $e_1$

dessert    $e_2$

is         $e_3$

excellent.  $e_4$

SoftMax

AVG $\longrightarrow \bigcirc$    MS.

受词出现次数的影响.

2): RNN、



SoftMax

$e_1$        $e_2$        $e_{10}$

Glove Word Embeddings.

$x_{ij}$: i 出现在 j 的上下文中的次数.

$$\text{minimize} \quad \sum_i \sum_j f(x_{ij})\left(\theta_i^T e_j + b_i + b_j - \lg x_{ij}\right)^2$$

加权: $x_{ij} = 0$; "the, a, an".

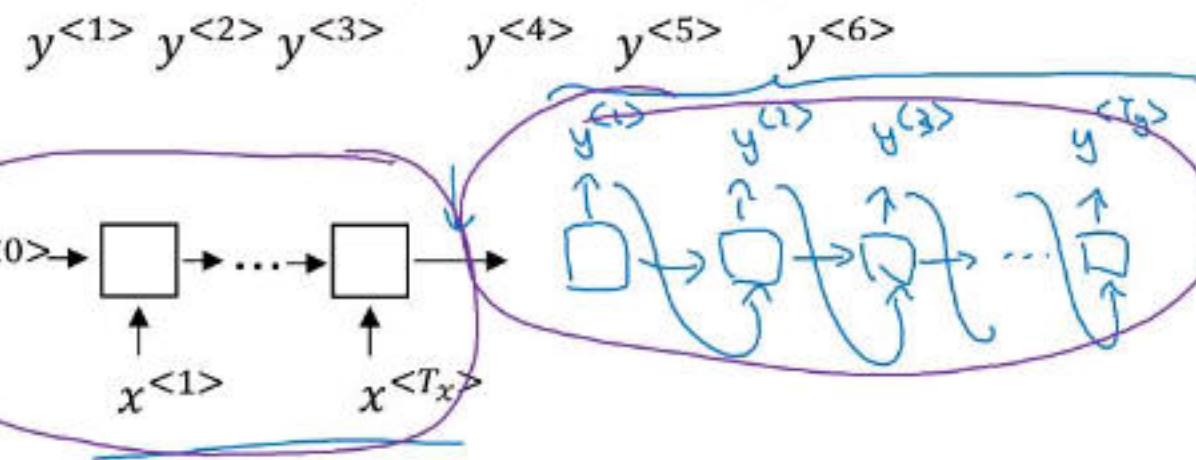$\theta_i$ 与 $e_j$ 是完全对称的, $\quad e_w = \dfrac{e_w + \theta_w}{2}$.

# Sequence to Sequence.

## Basic Models.

## Sequence to sequence model

$$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$$
Jane  visite  l'Afrique  en  septembre

$\longrightarrow$  Jane  is  visiting  Africa  in  September.

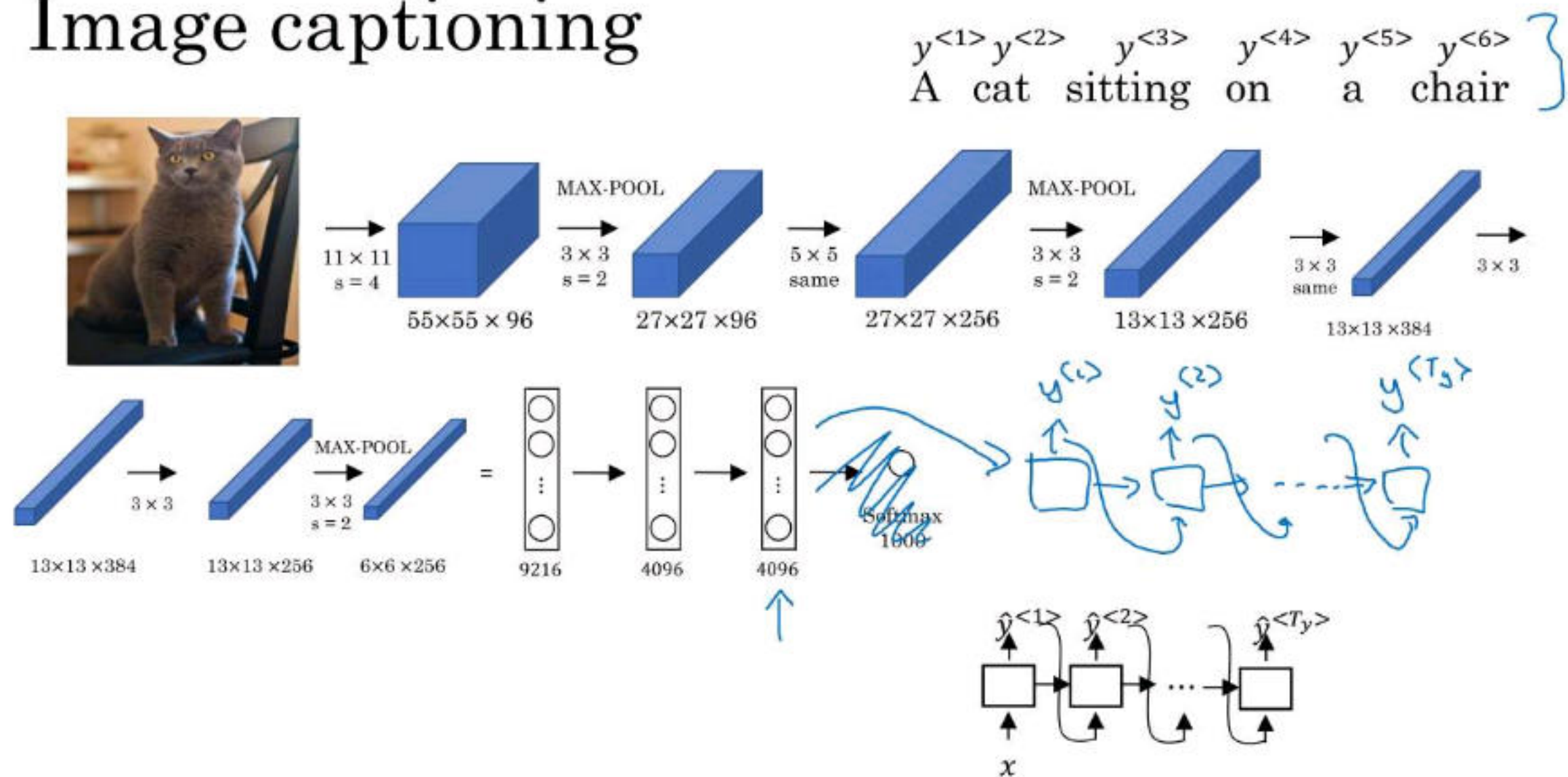$$y^{<1>} \; y^{<2>} \; y^{<3>} \qquad y^{<4>} \quad y^{<5>} \qquad y^{<6>}$$



[Sutskever et al., 2014. Sequence to sequence learning with neural networks]
[Cho et al., 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation]

Andrew Ng

Encoder — Decoder 结构 同样被用在了 Image Caption

## Image captioning

$$y^{<1>} y^{<2>} \qquad y^{<3>} \qquad y^{<4>} \quad y^{<5>} \quad y^{<6>}$$
A  cat  sitting  on    a    chair



[Mao et. al., 2014. Deep captioning with multimodal recurrent neural networks]
[Vinyals et. al., 2014. Show and tell: Neural image caption generator]
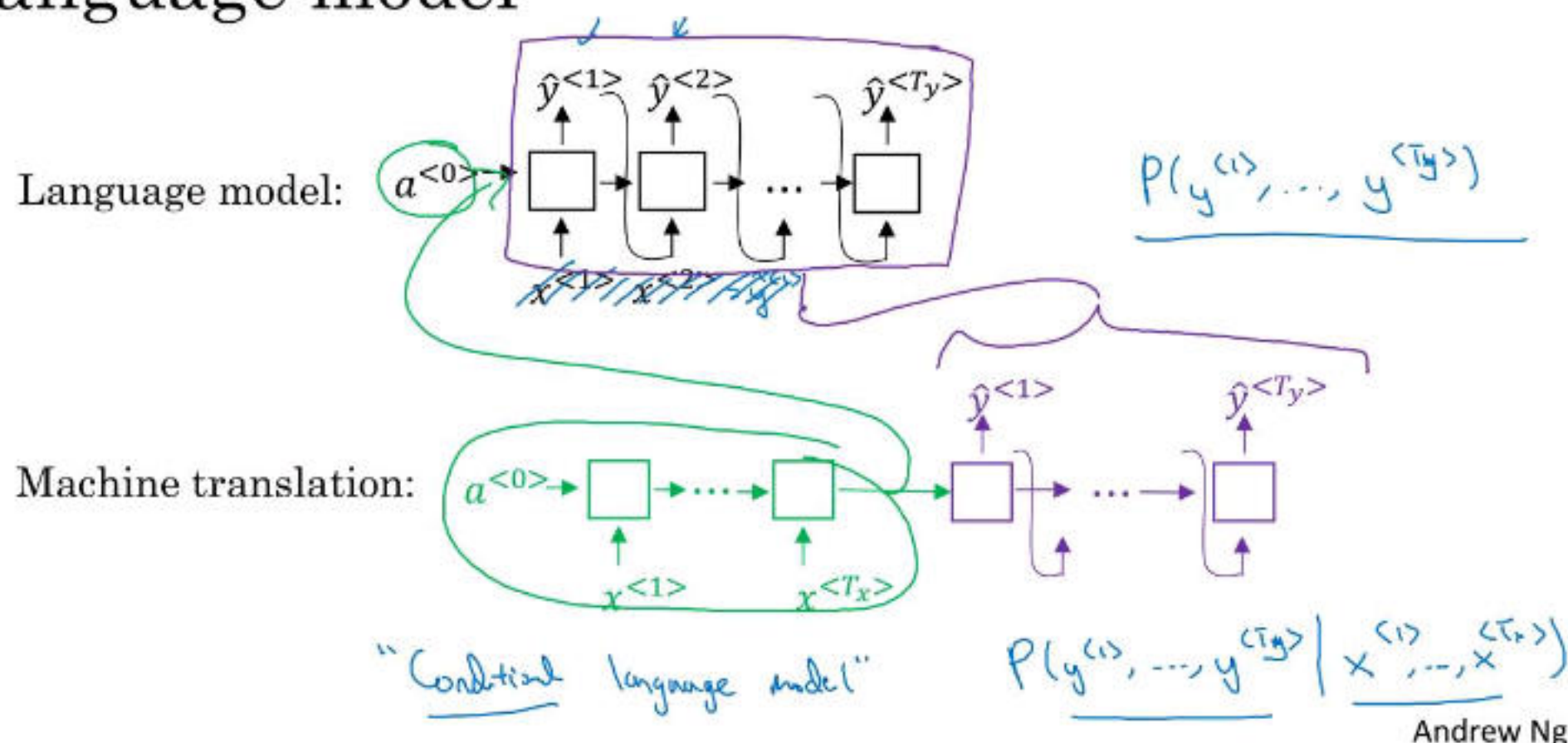[Karpathy and Li, 2015. Deep visual-semantic alignments for generating image descriptions]
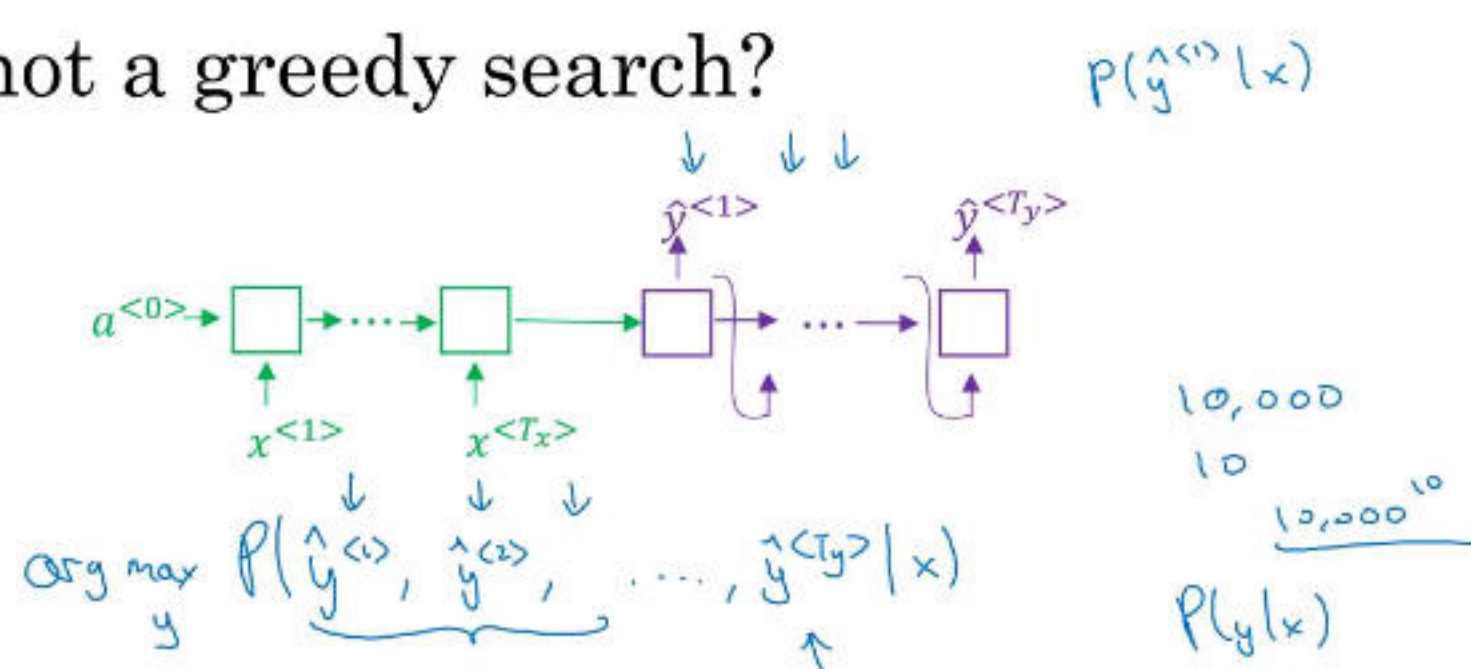
Andrew Ng

# Machine translation.

## Machine translation as building a conditional language model



Language model: $a^{<0>}$ ... $\hat{y}^{<1>}$ $\hat{y}^{<2>}$ $\hat{y}^{<T_y>}$

$x^{<1>},x^{<2>}...$

$P(y^{(1)},\ldots,y^{(T_y)})$

Machine translation: $a^{<0>}$ → ... → $x^{<1>}$ $x^{<T_x>}$ $\hat{y}^{<1>}$ ... $\hat{y}^{<T_y>}$

"Conditional language model"

$P(y^{(1)},\ldots,y^{(T_y)} \mid x^{(1)},\ldots,x^{(T_x)})$

条件语言模型. 源语句是条件.

## Why not a greedy search?

$P(\hat{y}^{(1)} \mid x)$



$a^{<0>}$ → ... → $x^{<1>}$ $x^{<T_x>}$ $\hat{y}^{<1>}$ ... $\hat{y}^{<T_y>}$

$\arg\max\limits_{y} P(\hat{y}^{(1)}, \hat{y}^{(2)}, \ldots, \hat{y}^{<T_y>} \mid x)$

10,000
10
$10,000^{10}$
$P(y \mid x)$

→ Jane is visiting Africa in September.
→ Jane is going to be visiting Africa in September.

$P(\text{Jane is goig} \mid x) > P(\text{Jane is visit} \mid x)$
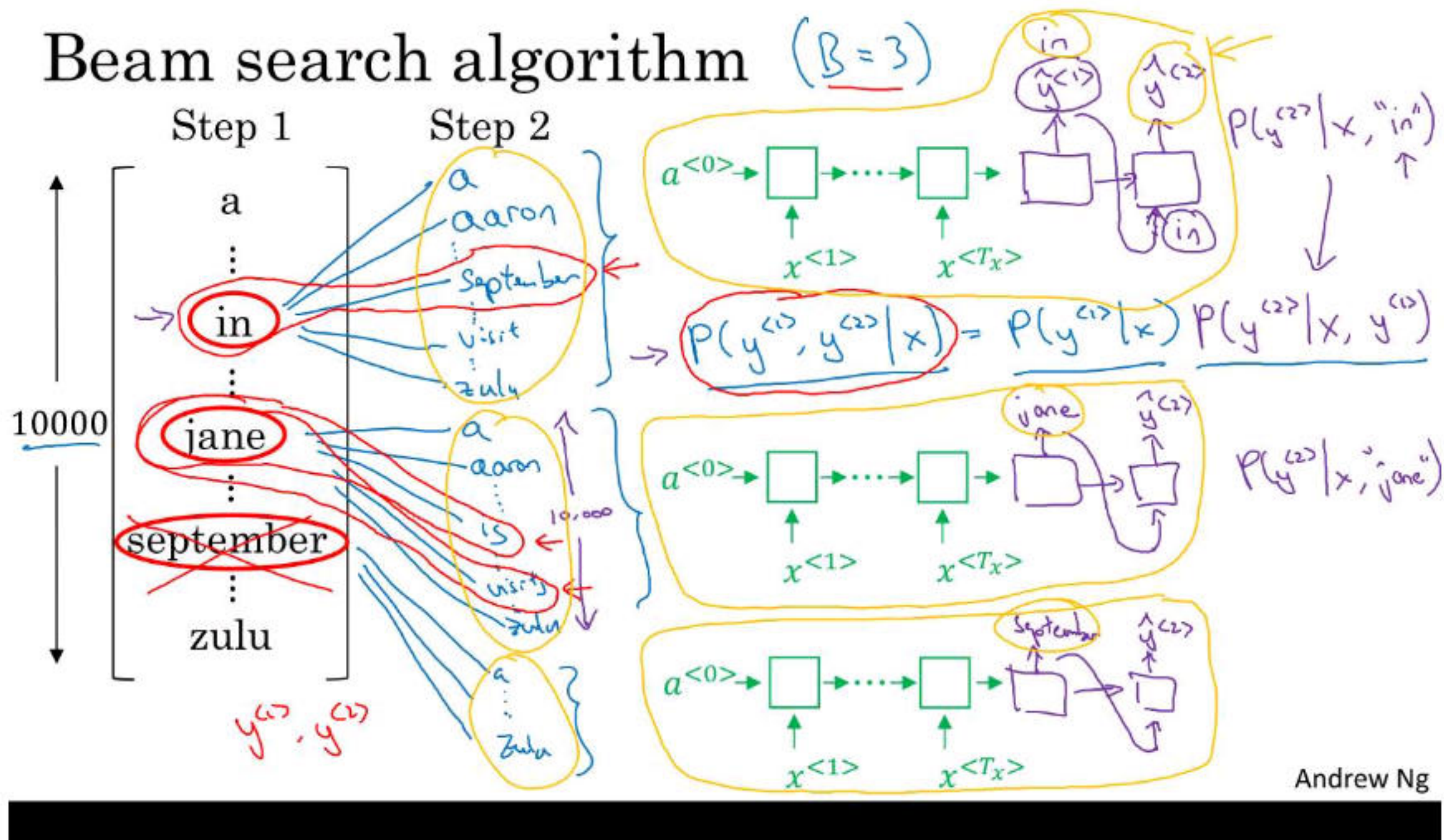
输出是随机采样的. $\arg\max\limits_{y} P(y \mid x)$
不能用 Greedy, 会非 全局最优.
↓
Beam Search.

Beam   Search.

每次保留   B = Beam   width 个候选词.



# Beam search algorithm $(B=3)$

Step 1     Step 2

Andrew Ng

RNN 给出了 $P(y^{<1>}|x)$     $P(y^{<2>}|x, y^{<1>})$

机器翻译想获得的是 $\arg\max P(\hat{y}^{<1>}, \hat{y}^{<2>})$

$P(y^{<1>}, y^{<2>}|x) = P(\hat{y}^{<2>}|\hat{y}^{<1>}, x) \cdot P(y^{<1>}|x)$

每次保留 B 个.

# Length normalization



$$P(y^{<1>}...y^{<T_y>}|x) = P(y^{<1>}|x) P(y^{<2>}|x, y^{<1>}) \cdots$$
$$P(y^{<T_y>}|x, y^{<1>}..., y^{<T_y-1>})$$

$$\arg\max_y \prod_{t=1}^{T_y} P(y^{<t>}|x, y^{<1>},...,y^{<t-1>})$$

$\log$

$$\arg\max_y \sum_{t=1}^{T_y} \log P(y^{<t>}|x, y^{<1>},...,y^{<t-1>}) \leftarrow$$

$\log P(y|x) \leftarrow$

$P(y|x) \leftarrow$

$T_y = 1, 2, 3, ..., 30.$

$$\frac{1}{T_y^{\alpha}} \sum_{t=1}^{T_y} \log P(y^{<t>}|x, y^{<1>},...,y^{<t-1>})$$

$\alpha = 0.7$     $\alpha = 1$
            $\alpha = 0$

Andrew Ng

对   $\arg\max P(y|x)$   正则 $\Rightarrow \arg\max \frac{1}{T_y^{\alpha}} \cdot \sum \log P(\cdot).$

## Error analysis on beam search

$P(y^*|x)$

$P(\hat{y}|x)$

Human: Jane visits Africa in September. ($y^*$)

Algorithm: Jane visited Africa last September. ($\hat{y}$)

Case 1: $P(y^*|x) > P(\hat{y}|x)$ ←

$\text{arg max } P(y|x)$
$y$

Beam search chose $\hat{y}$. But $y^*$ attains higher $\boxed{P(y|x)}$.

Conclusion: Beam search is at fault.

Case 2: $P(y^*|x) \lesssim P(\hat{y}|x)$ ←

$y^*$ is a better translation than $\hat{y}$. But RNN predicted $\boxed{P(y^*|x)} < \boxed{P(\hat{y}|x)}$.

Conclusion: RNN model is at fault.

Andrew Ng

对 Beam Search 出错的一些错误分析.

△ 机器翻译的评价指标: BLEU值.

Intuition: 与人类翻译越接近, 质量越好.

Bilingual Evaluating Understudy (替补).

需要截断. 得分, 否则 " the * ]".

n-grams:

$$P_n = \frac{\sum\limits_{\text{n-grams}\in\hat{y}} \text{Count-Clip(n-grams)}}{\sum\limits_{\text{n-grams}\in\hat{y}} \text{Count (n-grams)}}$$

$$BP = \begin{cases} 1 & , \text{if} \quad \text{output\_length} > \text{ref-length}. \\ \exp(1 - \text{output\_length}/\text{ref\_length}) & , \text{otherwise}. \end{cases}$$

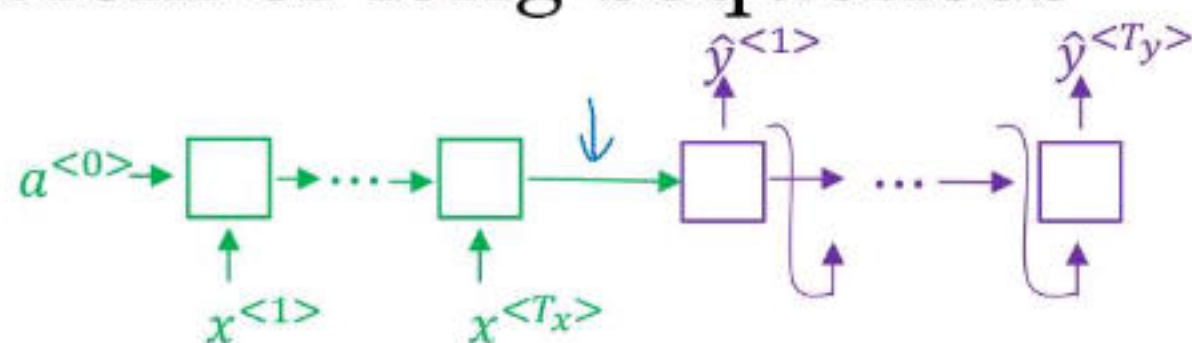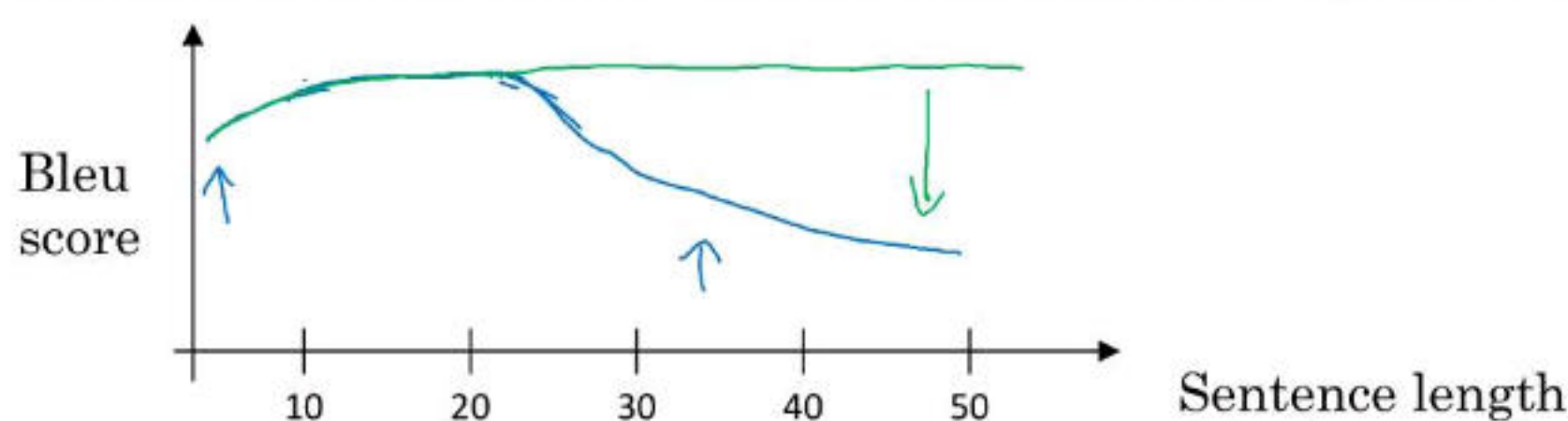$$BLEU = BP \cdot \exp\left(\frac{1}{K}\sum_{n=1}^{K} P_K\right)$$

brevity penalty.

用 exp 做了加权.

# Attention Mechanism.

## The problem of long sequences

$\hat{y}^{<1>}$    $\hat{y}^{<T_y>}$

$a^{<0>}$ → □ → ···· → □    ↓ → □ →  ···  → □

$x^{<1>}$        $x^{<T_x>}$

Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.
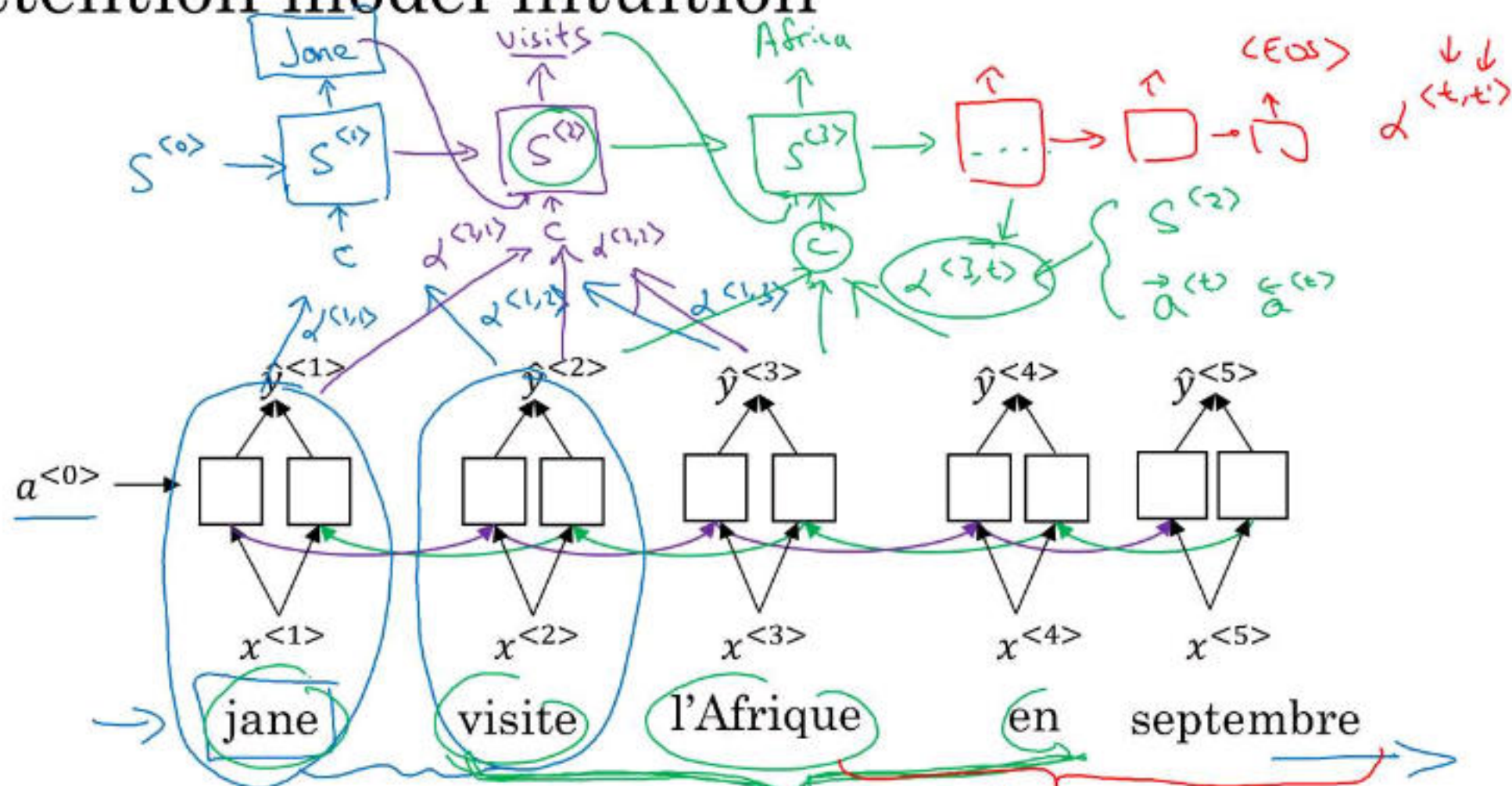
Jane went to Africa last September, and enjoyed the culture and met many wonderful people; she came back raving about how wonderful her trip was, and is tempting me to go too.

Bleu score

10    20    30    40    50    Sentence length

人做的是部分 翻译，而 RNN是 整句翻译（Bad）.

## Attention model intuition

Jane    Visits    Africa    <EOS>    $\alpha^{<t,t'>}$

$S^{<0>}$ → $S^{<1>}$ → $S^{<2>}$ → $S^{<3>}$ → □ ··· → □ → □

$S^{<2>}$

$\vec{a}^{<t>}$ $\overleftarrow{a}^{<t>}$

$\alpha^{<1,1>}$  $\alpha^{<1,2>}$  $\alpha^{<1,3>}$

$\alpha^{<1,1>}$  $\alpha^{<1,2>}$  $\alpha^{<1,3>}$    $\gamma^{<3,t>}$

$\hat{y}^{<1>}$    $\hat{y}^{<2>}$    $\hat{y}^{<3>}$    $\hat{y}^{<4>}$    $\hat{y}^{<5>}$

$a^{<0>}$ →

$x^{<1>}$    $x^{<2>}$    $x^{<3>}$    $x^{<4>}$    $x^{<5>}$
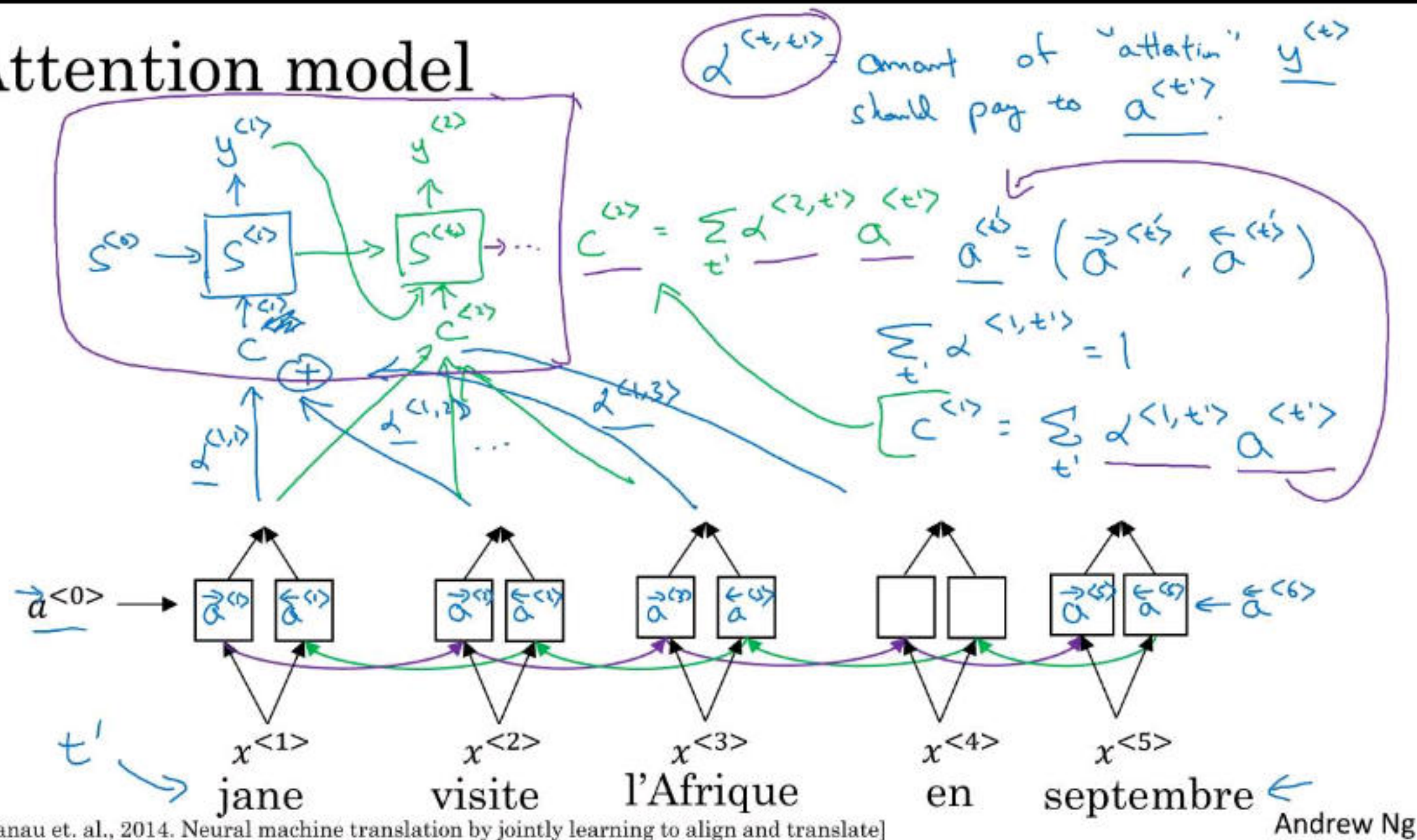
jane    visite    l'Afrique    en    septembre

[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

Att 模拟 部分 对部分的翻译.

# Attention model



$\alpha^{<t, t'>}$ — amount of "attention" $y^{<t>}$ should pay to $a^{<t'>}$.

$$c^{<2>} = \sum_{t'} \alpha^{<2,t'>} a^{<t'>}$$

$$a^{<t>} = (\vec{a}^{<t>}, \overleftarrow{a}^{<t>})$$

$$\sum_{t'} \alpha^{<1,t'>} = 1$$

$$c^{<1>} = \sum_{t'} \alpha^{<1,t'>} a^{<t'>}$$

$\vec{a}^{<0>} \rightarrow$

jane    visite    l'Afrique    en    septembre

$x^{<1>}$    $x^{<2>}$    $x^{<3>}$    $x^{<4>}$    $x^{<5>}$

$t'$

[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]    Andrew Ng

---

$$\sum_{t'}^{T_y} \alpha^{<t,t'>} = 1.$$

$$c^{<1>} = \sum_{t'} \alpha^{<1,t'>} a^{<t'>}.$$

$C$: Context.

$\alpha^{<t,t'>}$ : attention $y^{<t>}$ pay to $a^{<t'>}$.

---

# Computing attention $\alpha^{<t,t'>}$

$T_x \qquad T_y$

$\alpha^{<t,t'>}$ = amount of attention $y^{<t>}$ should pay to $a^{<t'>}$

$$\alpha^{<t,t'>} = \frac{\exp(e^{<t,t'>})}{\sum_{t'=1}^{T_x} \exp(e^{<t,t'>})}$$

$s^{<t-1>}$

$a^{<t'>}$

$\rightarrow e^{<t,t'>}$

$\alpha^{<t,t'>}$

$\hat{y}^{<t-1>}$    $\hat{y}^{<t>}$

$s^{<t-1>}$    $s^{<t>}$

$a^{<0>} \rightarrow$

$x^{<1>}$    $x^{<2>}$    $x^{<T_x-1>}$    $x^{<T_x>}$

[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]
[Xu et. al., 2015. Show, attend and tell: Neural image caption generation with visual attention]    Andrew Ng
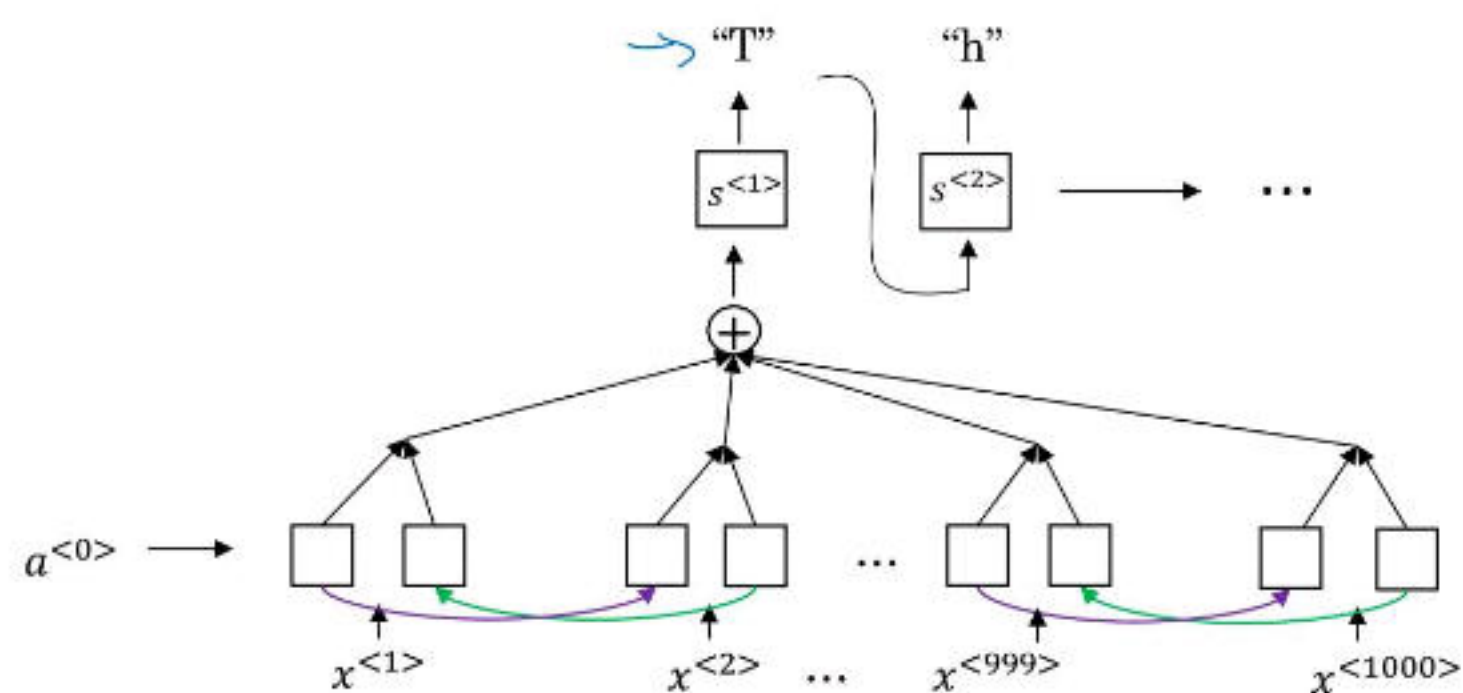
---

Att 会与 $s^{<t-1>}$, $a^{<t'>}$ 相关. 但不知. 构造小型 NN.

# Speech Recognition
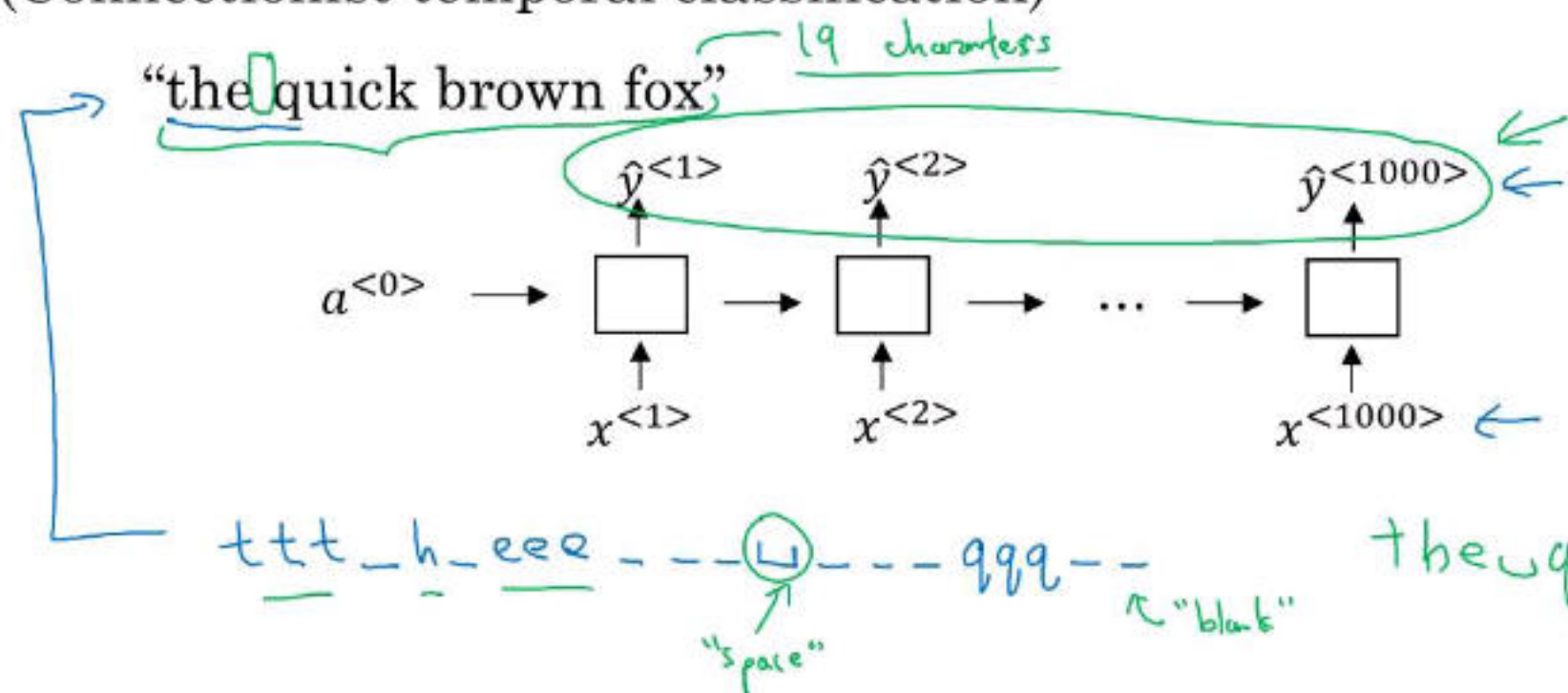
---

## Attention model for speech recognition



"T"    "h"

$s^{<1>}$    $s^{<2>}$ → ...

$a^{<0>}$ →

$x^{<1>}$    $x^{<2>}$ ...    $x^{<999>}$    $x^{<1000>}$

Andrew Ng

---

本质上也是一种翻译，也可以应用 Att。

---

## CTC cost for speech recognition

(Connectionist temporal classification)

19 charatess

"the quick brown fox"

$\hat{y}^{<1>}$    $\hat{y}^{<2>}$    $\hat{y}^{<1000>}$

$a^{<0>}$ → ☐ → ☐ → ... → ☐

$x^{<1>}$    $x^{<2>}$    $x^{<1000>}$

ttt_h_eee - - - - (ப) - - - qqq - -    theuq
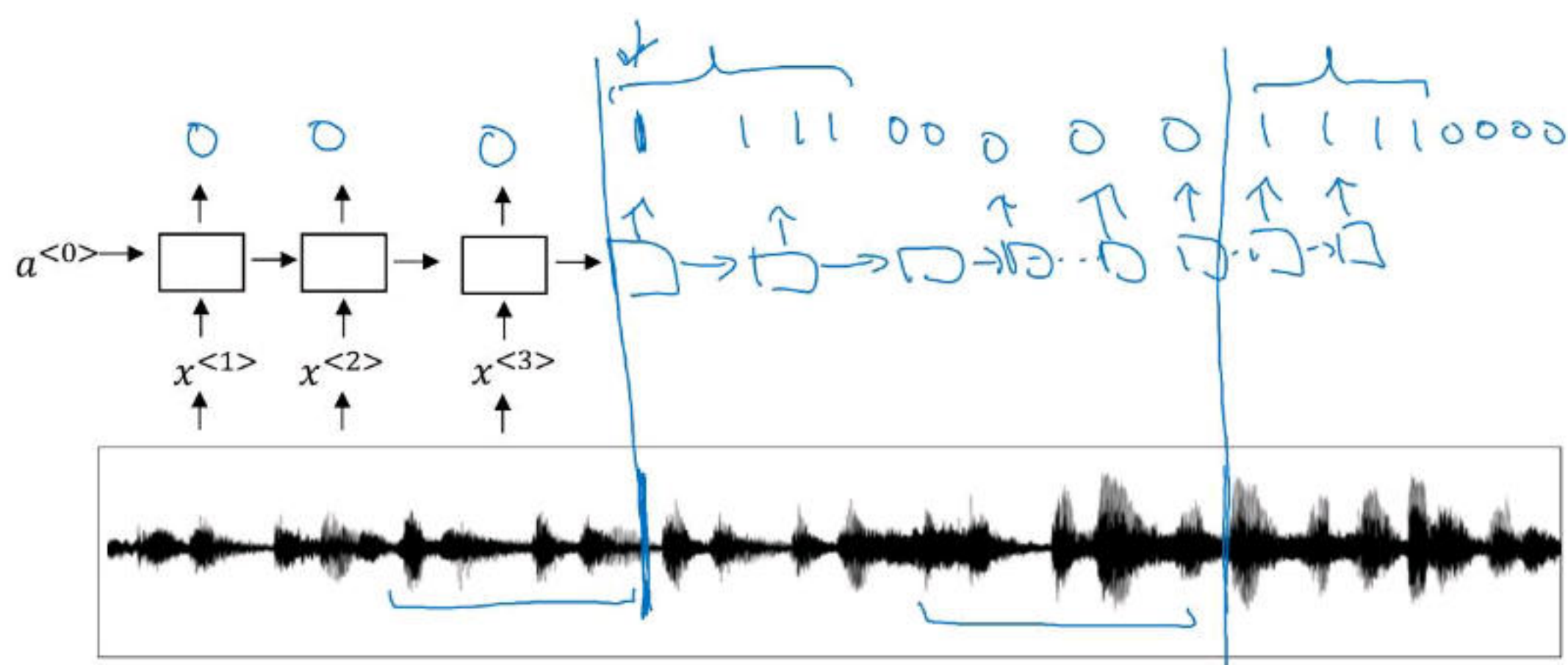
"space"    ↖ "blank"

Basic rule: collapse repeated characters not separated by "blank"

[Graves et al., 2006. Connectionist Temporal Classification: Labeling unsegmented sequence data with recurrent neural networks]    Andrew Ng

---

## Trigger word detection algorithm

用 0→1 标志
触发词频检测



O    O    O    1 11 0000 0 0 1 1 1 0000

$a^{<0>}$ →

$x^{<1>}$    $x^{<2>}$    $x^{<3>}$

Andrew Ng